

APPLICATION FOR UNITED STATES PATENT

MEMORY MANAGEMENT

By Inventors:

Gil Tene
1808 N. Shoreline Blvd.
Mountain View, CA 94043
A Citizen of the United States

Cliff N. Click
1808 N. Shoreline Blvd.
Mountain View, CA 94043
A citizen of the United States

Michael A. Wolf
1808 N. Shoreline Blvd.
Mountain View, CA 94043
A citizen of the United States

Ivan Posva
1808 N. Shoreline Blvd.
Mountain View, CA 94043
A citizen of Switzerland

Assignee: Azul Systems, Inc.

VAN PELT AND YI, LLP
10050 N. Foothill Blvd., Suite 200
Cupertino, CA 95014
Telephone (408) 973-2585

MEMORY MANAGEMENT

CROSS REFERENCE TO RELATED APPLICATIONS

This application claims priority to U.S. Provisional Patent Application No. 60/492,585 (Attorney Docket No. AZULP002+) entitled "STACK BASED
5 ALLOCATION" filed August 4, 2003 which is incorporated herein by reference for all purposes.

FIELD OF THE INVENTION

The present invention relates generally to computer programming techniques. More specifically, a technique for memory management is disclosed.

10

BACKGROUND OF THE INVENTION

Efficient memory management is important for efficient execution of computer programs. Many programming languages support memory management for heap memory. For example, some virtual machine implementations use a fast generational garbage collector that allocates objects from the heap, wherein the allocation streams
15 through memory. Typically, streamed data flushes lines out of the cache and forces extra read operations and write operations. This leads to extra memory traffic and inefficient utilization of caches.

One way of reducing memory traffic is to use stack based allocation for memory management. In systems that support stack based allocation schemes, objects are allocated on the running program's stack instead of a shared heap. The memory is freed when the stack frame exits. Since manual stack allocation requires extensive

5 programmer intervention and is error prone, stack based allocation is typically done automatically. Currently, successful stack based allocation requires accurate knowledge about the lifetime of objects. Such knowledge may be acquired using escape analysis. However, escape analysis is typically an expensive operation performed at compile time, making it impractical for large-scale development efforts. Also, current escape analysis
10 tends to be pessimistic, that is, it tends to limit stack based allocation to those objects that can be proven to be safe to allocate on the stack. Such pessimistic allocation potentially misses a significant number of effective stack based allocation opportunities.

It would be useful to have a memory management technique that would allow for efficient stack based allocation, without incurring the compile time overhead associated
15 with escape analysis. It would also be desirable if the memory management technique could maximize the opportunities for stack based allocation.

BRIEF DESCRIPTION OF THE DRAWINGS

Various embodiments of the invention are disclosed in the following detailed description and the accompanying drawings.

Figures 1A-1C are diagrams illustrating the stack in several stages.

5 Figure 2A is a diagram illustrating a stack that is modified according to one embodiment.

Figure 2B is a diagram illustrating another example of a stack that is modified upon detecting an escape.

10 Figure 3 is a flowchart illustrating a memory management technique using escape detection according to one embodiment.

Figure 4A is a diagram illustrating an address that incorporates a frame identifier according to one embodiment.

Figure 4B is a diagram illustrating address structures used in another embodiment.

15 Figure 5 is a flowchart illustrating escape detection and object reassignment according to one embodiment.

Figure 6 is a flowchart illustrating the recursive update of younger objects according to one embodiment.

Figure 7 is a flowchart illustrating how an object is moved to the reassignment frame according to one embodiment.

Figure 8 is a diagram illustrating the use of a display according to one embodiment.

5 Figure 9 is a flowchart illustrating an optimization process according to one embodiment.

DETAILED DESCRIPTION

The invention can be implemented in numerous ways, including as a process, an apparatus, a system, a composition of matter, a computer readable medium such as a computer readable storage medium or a computer network wherein program instructions
5 are sent over optical or electronic communication links. In this specification, these implementations, or any other form that the invention may take, may be referred to as techniques. In general, the order of the steps of disclosed processes may be altered within the scope of the invention.

A detailed description of one or more embodiments of the invention is provided
10 below along with accompanying figures that illustrate the principles of the invention. The invention is described in connection with such embodiments, but the invention is not limited to any embodiment. The scope of the invention is limited only by the claims and the invention encompasses numerous alternatives, modifications and equivalents. Numerous specific details are set forth in the following description in order to provide a
15 thorough understanding of the invention. These details are provided for the purpose of example and invention may be practiced according to the claims without some or all of these specific details. For the purpose of clarity, technical material that is known in the technical fields related to the invention has not been described in detail so that the invention is not unnecessarily obscured.

20 An improved technique for managing an object in memory is disclosed. In some embodiments, any attempt to place a reference to an object in an older frame is detected,

and the object is reassigned to a reassignment frame that is at least as old as the older frame. In some embodiments, the reassignment frame is the older frame. In some embodiments, a frame identifier is associated with the object and used for detecting any attempt to place a reference to an object in an older frame.

5 For the purposes of illustration, the following objects and related functions are used as an example:

```
10           Sam {  
              a = new;  
              bad = Uli();  
              Ted();  
              ...use bad...;  
          }  
  
15           Ted {  
              c = new;  
              return;  
          }  
  
20           Uli {  
              d = new;  
              return d;  
          }
```

25 The objects are allocated on the stack in this example. Figures 1A-1C are diagrams illustrating the stack in several stages. Figure 1A is a diagram illustrating the stack after Sam is invoked. A frame associated with Sam is allocated on the stack. The frame is a logical collection of memory used by the function Sam. Although the frame is shown as a contiguous block of memory in this example, it may sometimes be discontiguous and include several blocks.

Figure 1B is a diagram illustrating the stack after `bad = Uli()` is invoked. A frame is allocated for Uli after Sam's frame. In this example, a frame that is allocated later on the stack is deemed to be a younger frame than any frame allocated before it on the stack since the previously allocated frame will outlast the later allocated frame in terms of function scope. Thus, Uli's frame is younger than Sam's frame since Uli's frame will go out of scope before Sam's frame does.

In Uli's frame, an object "d" is allocated at the beginning of Uli's frame. A pointer in Sam, "bad", points to "d". Assigning Uli to pointer "bad" has the effect of assigning object "d" of Uli to a pointer in Sam; this is referred to as an escape. An escape happens whenever an object allocated in a younger frame becomes visible to an older frame.

If not properly handled, escapes can lead to undesirable effects, as illustrated in Figure 1C. Figure 1C is a diagram illustrating the stack after `Ted()` is invoked. Since Sam's frame is an older frame that outlasts Uli's frame, once the call to Uli completes, Uli's frame is cleared from the stack. Ted's frame is then allocated on the stack at the beginning from the same location as where Uli's frame was. However, the pointer in Sam's frame, "bad", still points to the location where the "d" object in Uli's frame once occupied. The address occupied previously by the "d" object in Uli's frame is now occupied by the "c" object in Ted's frame, thus when `bad` is invoked, unpredictable behavior may occur.

In some embodiments, escape detection is performed to prevent unpredictable behavior associated with escapes. Once an escape is detected, it may be handled by modifying the stack frames and objects to ensure proper program code behavior. Figure 2A is a diagram illustrating a stack that is modified according to one embodiment.

- 5 Returning to the previous example, when an escape of Uli's frame is detected (i.e., when `bad = Uli()` is invoked), Sam's frame is modified to store the object "d" directly. In this example, "d" is copied onto Sam's frame and the pointer "bad" now points to the address in Sam's frame that is now occupied by "d". When Uli's frame goes out of scope, the pointer to "d" is still valid since "d" has become a part of Sam's frame and shares the
- 10 same scope as other objects in Sam's frame.

- Figure 2B is a diagram illustrating another example of a stack that is modified upon detecting an escape. In this example, the size of object "d" is determined to be greater than the amount of space available in Sam's frame for storing additional objects. Thus, an additional block of memory is allocated in an overflow area to store "d", and
- 15 "bad" now references the newly allocated "d". Sam's frame is a logical frame that includes the memory originally allocated for the frame and any additional memory used to store modified objects. In some embodiments, references to "d" and references to younger objects by "d" are updated also. Details of the updates are discussed below.

- Figure 3 is a flowchart illustrating a memory management technique using escape
- 20 detection according to one embodiment. When an object is created, it is assigned to an initial frame (300). Any attempt to place a reference to the object in an older frame is detected (302). The object is then reassigned to a reassignment frame that is at least as

old as the older frame. In some embodiments, the reassignment frame is the older frame; in some embodiments, the reassignment frame is a frame older than the older frame.

In some embodiments, the system learns upon detecting an escape and uses the reassignment information to improve performance. For example, the system may learn
5 that an escape occurs at a certain allocation site where the object is allocated. As used herein, an object's allocation site refers to the location in the program code that allocates the object. To reduce the cost of reassigning the object next time, program code at the allocation site of the object is modified so that the object is automatically reassigned next time it is allocated (306). In some embodiments, the same type of object may have
10 multiple allocation sites that are distinguished according to their call paths. The call path information can be stored and code at each allocation site can be modified differently to handle escapes.

In some embodiments, the system learns when additional memory is required for the reassignment frame to store the reassigned object. The frame creation site (i.e. the
15 location in the program code used to create frames) of the reassignment frame may be modified so that next time the frame is created, additional memory is automatically allocated for the frame to accommodate the reassigned object (308). In some embodiments, side tables or other appropriate data structures may also be used to assist escape handling and object creation. These learning features are optional and may be
20 omitted.

Frame identifiers (FIDs) are sometimes used in conjunction with references or address of the objects to facilitate escape detection. Figure 4A is a diagram illustrating an address that incorporates a frame identifier according to one embodiment. In this example, several bits at the beginning of an address block are reserved to store the FID.

5 The rest of the bits are used to store the actual address or reference to an object. The value of the FID is a function of the object's frame depth in some embodiments. For example, in an embodiment where the function is a count of the frame depth, a program with 5 nested calls results in a FID of 5 for the newest frame. Escape attempts, including any attempt to store an object allocated in a younger frame to another object allocated in
10 an older frame as well as returning an object allocated in a younger frame to a method invoked from an older frame can be detected by comparing the FIDs.

In some embodiments, heap memory is regarded as a regular frame and assigned a FID; in some embodiments, heap memory is treated specially. Figure 4B is a diagram illustrating address structures used in another embodiment. In this example, a heap
15 identifier is used to differentiate objects allocated on the heap from the object allocated on the stack. The heap identifier may be a distinguishing bit or a distinguishing set of bits. In this example, the heap identifier is a single reserved bit. An address for an object allocated on the heap such as 400 sets the heap identifier to 1 and the rest of the bits are used to store the address or reference to the heap memory. An address for an object
20 stored on the stack such as 402 sets the heap identifier to zero indicating that a predetermined number of bits following the heap identifier are used as a FID and the rest

of the bits are used to store the stack address or reference to the object. Using a heap identifier makes more bits available for storing the heap address.

Another way to detecting an escape is to directly compare the addresses of the objects. In some embodiments, the direction in which the stack grows determines how the comparison is to be performed. For example, if the stack grows up, an object with a higher address value is deemed to be younger than an object with a lower address value; if the stack grows down, the object with a higher address value is considered to be older than an object with a lower address value. Additional checking may be performed to detect the false positive created by objects that are allocated in the same frame, but in the reverse order. In some embodiments, a FID is stored just prior to the object with a negative offset in the stack or as part of the object's header word, and the additional test compares the FIDs to determine the relative age of objects.

Figure 5 is a flowchart illustrating escape detection and object reassignment according to one embodiment. In this example, the FID of the object is compared with the FID of the frame that references the object (500). It is determined whether the reference frame FID is older than the object FID (502). If the reference frame FID is not older, there is no escape and no action needs to be taken (512). If it is determined that the reference frame FID is older than the object FID, an escape has occurred and the object is moved to its reassignment frame, in this case the reference frame (504). The object FID is then reset to the FID of the reassignment frame (506). The references to the object are also updated (508). In this embodiment, since the object has not yet escaped, the references to be updated include the ones in the reassignment frame and any frames

younger than the reassignment frame. Younger objects referenced by the reassigned object are then updated recursively (510).

In some embodiments, the escape detection is performed with assistance from hardware configured to make the detection more efficient. For example, hardware assist
5 may be used to compare the FID's of the object and the reassignment frame, and generate a trap if an escape is detected. Such hardware can be specially designed and/or configured to efficiently perform operations such as FID comparison and increase system speed.

Figure 6 is a flowchart illustrating the recursive update of younger objects
10 according to one embodiment. In this example, it is determined if the next younger referenced object is available (600). If no more referenced object is available, the recursive call returns (602). If a younger object is available, it is selected (604). The escape detection and object reassignment process shown in Figure 5 is then repeated for the selected referenced object and the recursion continues (606).

15 Figure 7 is a flowchart illustrating how an object is moved to the reassignment frame according to one embodiment. It is determined whether the object fits in the reassignment frame (700). If the object fits in the reassignment frame, it is copied into the reassignment frame (710), and the used space is traced so that it will not be reassigned to other objects (712). If the object does not fit, an overflow area for the
20 frame is created (702). The object is then copied to the overflow area (704). In this example, one overflow area per frame is created. If a frame stays in scope for an

extended period of time, it may accumulate multiple objects in its overflow area. Objects that are no longer used should be removed. In some embodiments, it is determined whether the number of overflow objects has exceeded a certain threshold (706). If the overflow has not exceeded the threshold, there is still space available and no action needs
5 to be taken (714). If, however, it is determined that there are too many overflow objects, the thread-local garbage collector then removes the dead objects from the overflow area (708). In some embodiments, instead of using an overflow area, the reassignment frame is expanded to accommodate the object.

In some embodiments, a display is used to make the process of reassigning the
10 object more efficient. Figure 8 is a diagram illustrating the use of a display according to one embodiment. In the example shown, several objects, main, Ted, Mary and Bar are shown in different frames on stack 802. If, for example, Bar escapes to Ted's frame, determining where to copy Bar can incur significant operational overhead since the objects and the frames are of different sizes. In this embodiment, a display 800 is used to
15 assist this process. Display 800 is an array that stores information about the frames and objects. In some embodiments, the array entries include the addresses for each of the frames. In some embodiments, the array entries include the addresses for the location in each object where an escaped object can be copied. Since the array element size is fixed, finding the location to copy the escaped object now involves going to a fixed offset in the
20 display array to look up the address.

In some embodiments, the program code may be optimized using escape information. Figure 9 is a flowchart illustrating an optimization process according to one

embodiment. In this example, the process begins with running code that has not been optimized, such as the interpreted code (900). Escape data is then gathered (902). In some embodiments, for each allocation site, the maximum number of frame levels any object created at the allocation site escapes (also referred to as the maximum escape
5 depth) is recorded. This may be achieved by, for example, tagging each object with its allocation site information and associating an escape counter with each allocation site to store information about the maximum escape depth. When an object escapes, its allocation site is found and the escape counter associated with the allocation site is updated if the object's escape levels exceed the maximum escape depth currently stored
10 in the escape counter of the allocation site.

Data regarding “hot code” is also gathered (904). A section of code is determined to be “hot” if it is frequently used. The criteria for determining which code section is hot may vary for different embodiments. For example, a counter may be used to determine how frequently a code path is exercised. A fast compilation is then performed to generate
15 a version of compiled code that is optimized a bit more than the original interpreted version (906). In some embodiments, the compiler performs a small amount of inlining, including inlining hot code based on the information previously gathered. For example, if an object is included in a section of hot code, and if the escape depth of the object is below a certain threshold, then the call context of the object up to the escape depth may
20 be inlined. In another example, if the escape data indicates that an object escapes all the way to the heap, then the compiler may choose not to inline calls to the object to save memory.

The code obtained from the fast compilation is then run (908), and the escape data on the inlined code is gathered (910). A slower optimizing compilation is then performed using the escape data gathered to provide a more optimized program (912). In some embodiments, the compiler provides deeper inlining for the hot code. Note that in order
5 for the deeper inlining not to fail, objects called in the inlined code should not escape further than the context that is inlined. Although escape detection detects escapes that occur in the parts of the program code that have been exercised, it may not provide escape information for the rest of the code that has not been run. An allocation site may create objects that escape more levels than the current maximum escape depth, thus
10 inlining the code that creates the object only up to the current maximum escape depth may result in failure. In some embodiments, a more accurate escape assessment is obtained by performing an escape analysis. The escape analysis, although expensive to run, can accurately determine whether objects created by each allocation site escapes further than the current prediction. In some embodiments, the compiler uses the results
15 of the escape analysis to further optimize the code. For example, if the objects do not escape further than the current maximum escape depth, the code associated with the object can be safely inlined everywhere; if the object escapes further, deeper inlining may be performed.

In some embodiments, rather than using the same maximum escape depth to
20 inline all instances of the same type of object to the same depth, multiple allocation sites for the object are tagged and different escape depths are used to more efficiently inline code that calls the object. In one example, an instance of a function named Bar is called

by Sam and Bar creates objects that escape one frame level up to Sam. Another instance of the function Bar is called by Mary, which is in turn called by Ted, and here Bar creates objects that escape two frame levels up to Ted. Rather than inlining all calls to Bar at least two frame levels deep, the allocation site of Bar associated with Sam is inlined one
5 frame level deep and the allocation site of Bar associated with Mary is inlined two frame levels deep into Ted.

Although the foregoing embodiments have been described in some detail for purposes of clarity of understanding, the invention is not limited to the details provided. There are many alternative ways of implementing the invention. The disclosed
10 embodiments are illustrative and not restrictive.

WHAT IS CLAIMED IS: